



**QUEEN'S
UNIVERSITY
BELFAST**

Local Anomaly Detection by Application of Regression Analysis on PMU Data

Rafferty, M., Brogan, P., Hastings, J., Lavery, D., Liu, X., & Khan, R. (2018). Local Anomaly Detection by Application of Regression Analysis on PMU Data. In *Proceedings of Power and Energy Society General Meeting (PESGM), 2018* (IEEE Power & Energy Society General Meeting: Proceedings). Institute of Electrical and Electronics Engineers Inc..

Published in:

Proceedings of Power and Energy Society General Meeting (PESGM), 2018

Document Version:

Peer reviewed version

Queen's University Belfast - Research Portal:

[Link to publication record in Queen's University Belfast Research Portal](#)

Publisher rights

© 2018 IEEE.

This work is made available online in accordance with the publisher's policies. Please refer to any applicable terms of use of the publisher.

General rights

Copyright for the publications made accessible via the Queen's University Belfast Research Portal is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The Research Portal is Queen's institutional repository that provides access to Queen's research output. Every effort has been made to ensure that content in the Research Portal does not infringe any person's rights, or applicable UK laws. If you discover content in the Research Portal that you believe breaches copyright or violates any law, please contact openaccess@qub.ac.uk.

Local Anomaly Detection by Application of Regression Analysis on PMU Data

Mark Rafferty, Paul Brogan, John Hastings, David Lavery, Xueqin (Amy) Liu, Rafiullah Khan
School of Electronics, Electrical Engineering and Computer Science, Queen's University Belfast, UK

mrafferty16@qub.ac.uk

Abstract—PMU data has the potential of providing a wealth of information on power system operation, health, faults and anomalies. PMU tend to provide tens of measurements per second, therefore automated anomaly detection is required; especially for use in real or near-time applications by power system operators. This paper demonstrates a method of detecting local anomalies in PMU data utilizing multiple linear regression. A window of near-time data is employed to generate a regression function that predicts the live data that arrives. If the error between the observed and predicted values exceeds a threshold an exception is noted. The threshold is dynamically updated based on the error in the regression function, allowing the method to work equally well on data of varying regularity. This anomaly detection method is not tuned to particular events and should detect novel occurrences. The method is evaluated on two numerical case studies, a genuine power system event and a man in the middle cyber attack. Real data was collected from PMUs placed on the Irish power system.

Index Terms—Anomaly detection, PMU, machine learning, linear regression, cyber security

I. INTRODUCTION

Many methods exist for event detection and these are often employed in digital fault recorders (DFR). A typical method is to set simple thresholds for frequency, RoCoF, sequence components, harmonics or magnitude variations. Care must be taken when choosing thresholds as the region between missing events and near constant triggering can be small. Although typical thresholds will vary between power systems, they will also vary between sites, meaning each placement may need to be separately evaluated and monitored; preventing large scale rapid deployment. The method demonstrated in this paper is adaptive both to the local environment, both spatially and temporally, as the event threshold is updated based on the data recorded in the previous seconds.

As well as legitimate power system events, other anomalies can appear in the synchrophasor information due to 'bad data.' This can happen in a number of ways. For example, a missing GPS signal can cause deviation in phase-angle estimation due to a misaligned phase-locked loop. Information on GPS lock is usually sent in the PMU data stream, but unless this is processed by the application, the error can persist through the phasor estimation stage. Hardware or software related data acquisition errors are also possible due to hardware failure, or where mis-configured analog-to-digital (ADC) modules or software processing (pre-phasor estimation) lead to errors in the output of the phasor-estimation section of the PMU.

Another, more sinister possibility is that the PMU data had been modified in transit deliberately. The most prevalent standard in use today for synchrophasor transport is IEEE C37.118.2 (2005 or 2011) [1]. This PMU data transport protocol contains no security mechanisms at all, [2], and is highly susceptible to 'in-transit' manipulation. Considering the historic use of PMU data for post-fault analysis, this may not seem like much of an issue. However, PMU has the potential to become widely used for real-time analysis and as inputs to advanced control systems - therefore this type of attack vector requires consideration and system design.

Whether it's a legitimate event, an unlikely data acquisition-level fault, or an even less likely targeted data-stream manipulation, the method presented in this paper detects anomalous activity that may cause inappropriate PMU response. It will allow anomalous data to be quickly 'flagged' for further analysis.

To achieve high accuracy in the detection of anomalies in PMU data, this paper has (1) employed Multiple Linear Regression Analysis for the prediction of a power system variable based on previous measurements and other recorded variables in the power system; (2) developed a metric for the threshold of detecting anomalies in PMU data based on the difference between observed and predicted values; (3) implemented a sliding window approach to the methodology to dynamically update thresholds for anomaly detection using current, relevant power system parameters; (4) validated the proposed method through two numerical case studies of anomalous data collected from the Irish power system, a genuine power system event and a MITM attack on system frequency measurement on PMU data.

II. MACHINE LEARNING FOR ANOMALY DETECTION

Machine learning [3] is described as the construction of methods which automatically improve with experience. These can be split into two main types, supervised and unsupervised. Supervised learning is when the training data consists of examples of inputs and corresponding targets (typical problems are classification and regression), whilst unsupervised learning is when the training data consists of inputs only with targets unknown but worked out by the algorithm (clustering is a typical problem). A typical machine learning problem consists of two main phases, these are learning and prediction.

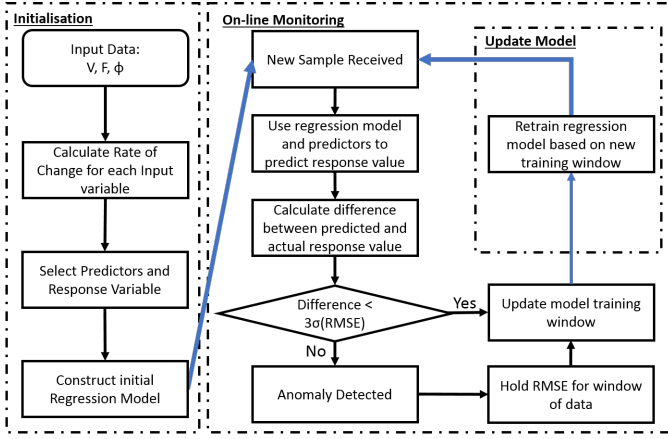


Fig. 1: Flow Chart for the proposed methodology for anomaly detection on PMU data using multiple regression analysis

A. Regression Analysis

Regression Analysis [4] is a technique in machine learning that is used to determine the relationship between two (simple regression) or more (multiple regression) variables. The goal of regression analysis is to determine the value of parameters for a function that gives the best fit line through a set of observations (learning phase), thus allowing the prediction of one variable based on the other recorded variables and the regression model (prediction phase). In this section the two most basic and commonly used regression types will be presented, these are simple and multiple linear regression.

A.1. Simple Linear Regression

In this case the model is bivariate, and shows the relationship between one independent variable (predictor), x , and a dependent variable (response), Y , is given by:

$$Y = \beta_0 + \beta_1 x + \epsilon \quad (1)$$

β_0 and β_1 represent the model coefficients, with β_0 representing the intercept and β_1 the coefficient of x (or the parameter of the slope), and ϵ the error in the model. A regression model can then be determined by learning the values of β_0 and β_1 from equation 1 for a training dataset of x and y values. Using the regression model a new Y sample can be predicted using the corresponding x value for the sample and equation 1.

A.2. Multiple Linear Regression

Multiple predictors are used in the model to build the relationship between predictors and response variables. Denoting a set of samples in a dataset at the i -th sample instant as $\mathbf{z}_i \in \mathbb{R}^{1 \times n}$, a data matrix $\mathbf{Z} \in \mathbb{R}^{i \times n}$ can be constructed with each row representing a sample, where i is the number of samples and n the number of variables in the matrix \mathbf{Z} . If a single variable is selected as a response (y) and the other variables selected as predictors ($x_1, x_2, x_3, \dots, x_{n-1}$), then for the i -th sample the regression equation is given as:

$$y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \beta_3 x_{i,3} + \dots + \beta_{n-1} x_{i,n-1} + \epsilon_i \quad (2)$$

where, similarly as before, β values represent the model coefficients, with β_0 representing the intercept, and $\beta_1 \dots \beta_{n-1}$ representing the coefficient of each element of \mathbf{X} , where

$\mathbf{X} = x_{i,1} \dots x_{i,n-1}$ respectively, and ϵ the error. Again, once the values of β have been calculated for a training dataset of \mathbf{X} and y values, a regression model can then be used for the prediction of y values given \mathbf{X} .

B. Adaptive PMU Anomaly Detection

The basic process for the proposed method is presented in Fig 1, and is separated into three main components. These are: 1) Initialization; 2) On-line Monitoring; and 3) Update Model. The initial stage of the method allows the construction of the initial regression model, and requires power system data for frequency (f), phase angle (ϕ) and voltage (v) and their rate of change (RoC) is calculated for each variable (Δf , $\Delta \phi$ and Δv respectively).

In the study presented the variables that were chosen for the construction of the initial regression model were:

Predictors

- Frequency
- Rate of change of frequency (RoCoF)
- Voltage
- Rate of change of voltage (RoCoV)

Response

- Rate of change of phase (RoCoP)

After the initial regression model has been constructed, it can be used to predict the next response value from the predictor values.

B.1. Sliding Window Regression

An adaptive regression technique was required to allow the model to be re-trained as power system conditions are continuously changing. The methodology incorporates a sliding window [5] approach. The model is continuously retrained as new data is received. Important parameters for optimization are window size and how often model retraining occurs.

B.2. Anomaly Detection - 68-95-99.7 rule

In literature the method of detecting anomalous data by calculating standard deviation from the mean has been extensively researched [6], and can be known as the 68-95-99.7 rule. In normally distributed datasets the amount of data that lies between 1 standard deviation of the mean is 68%, 2 standard deviations is 95% and 3 standard deviations is 99.7%.

The standard deviation is calculated from the RMS error, between the observed and predicted data within the sliding window, and is expressed as:

$$\sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}} \quad (3)$$

where y_i and \hat{y}_i are the observed and predicted values for the i -th observation and n the number of samples.

The choice of the number of standard deviation from the mean of RMS error was investigated, and findings are presented in Fig 4. This figure displays an initial regression model for normal power system operating data, with the top plot showing the observed and predicted RoCoP values whilst the bottom shows the difference and also a number of standard



Fig. 2: Location of Queen’s University Belfast PMUs on Irish power system

deviation thresholds. A trade off between a threshold which is not prone to nuisance tripping (false alarms) and has a small non detection zone (NDZ) is required to detect when an actual anomaly has occurred in the PMU data. From Fig 4, 3 standard deviations was selected as the threshold.

To reduce potential false alarms, and also outliers in the data causing a trigger, a time delay of 0.5 seconds is implemented into the process, meaning that the threshold must be violated for more than 0.5 seconds before a genuine anomaly is detected. Once a genuine anomaly has been detected, the model is held constant at pre-anomaly values, until the anomaly has been cleared to stop contamination of the regression model with anomalous data.

III. EVALUATION WITH REAL POWER SYSTEM DATA

To illustrate the potential of the proposed methodology for anomaly detection in local PMU data, it’s performance was evaluated for real power system data collected from the Irish Power system. These PMUs are a combination of commercial and open source, developed at Queen’s University Belfast [7], [8], and the location of each is highlighted in Fig 2.

GPS timing errors, and faulty hardware should not reach the output of the PMU if the system is designed correctly so will remain outside the scope of testing for now. However, the data-stream manipulation does present a very real potential attack vector in the industry, particularly if the attacked PMU streams are used in any control or real-time ‘cyber-physical’ systems. If PMUs are deployed in a large-scale - particularly with an increase in distributed energy resources (DERs), it’s very likely that some of these data streams will require traversal over less secure networks such as the internet. If these streams are not hardened, then the current PMU communications protocol is wide-open to manipulation.

To demonstrate this, a wide-area PMU transport simulator has been developed in the laboratory which allows wide-area communications to be accurately simulated. The test-bed

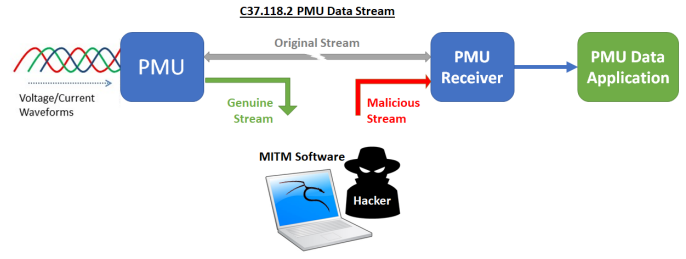


Fig. 3: Overview of MITM attack on PMU Data

comprises of a number of PMUs, a commercial PMU and an OpenPMU. The system allows for data in either of these streams to be modified in-transit, with frequency, phase angle and voltage magnitude modifications of the data possible.

The attack works by what is known as a MAC spoofing Man-In-The-Middle (MITM) attack (Fig 3). Media Access Control (MAC) is used at the data-link layer in computer communications, it is the protocol that allows communication in a switched computer network environment. The attacker essentially spoofs the MAC address (unique to each computer in a Local Area Network) of the PMU and PMU data receiver, becoming a ‘middle-man’ in the stream, without either party knowing. The ‘attacker’ is essentially a Python script running on a Linux machine. The unencrypted/unsigned nature of the most common PMU data format makes it very susceptible to this type of attack (IEC 61850-90-5 addresses these issues, as detailed in [2], but is not commonly deployed yet).

Spoofing attacks like this require access to the local network at either the PMU sending or receiving ends. This may sound difficult to achieve, given that the locations of such devices are often off-limits to anyone but authorized personnel. However, access to LAN environments can often be achieved by attackers through back-doors opened through other means - perhaps by malware infection following a click on a ‘phishing’ email. If unsafe data streams are to remain in use for the short-medium term, then the proposed anomaly detection method can be used to ‘flag’ messages for further analysis that aren’t simply power systems phenomena.

For this study a window size of 4 seconds (equating to 40 frames for PMUs reporting at 10 Hz) was used with retraining of the regression model occurring every 1 second (or 10 frames). Meaning that the previous 4 seconds of data is used to construct the regression model, which is then used, in conjunction with the predictor variable values, to predict the next 1 second of response data before retraining occurs. This process is then continuously repeated.

A. Construction of Initial Regression Model

To begin with historical, non-anomaly data is required for construction of the initial regression model. From this a window of data is selected and separated into their respective predictors and response data, with β values calculated for each predictor and the intercept. Once the regression model has been obtained, it is natural to evaluate how effective it is at summarizing the relationship between \mathbf{X} and y , and thus how good it will be at predicting y by using the RMS error.

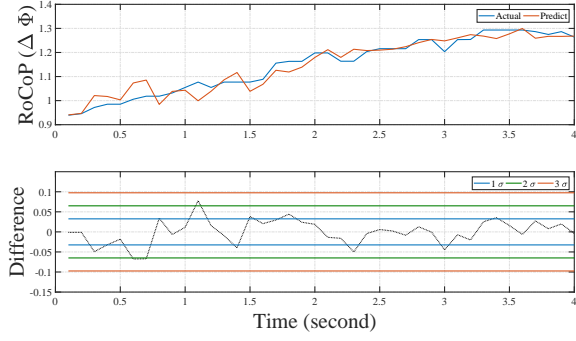


Fig. 4: Investigation into threshold settings for anomaly detection for normal power system operation data

The dataset used for model training is used to see how effective it is in the prediction of itself. This is illustrated in Fig 4, showing observed and predicted response values of RoCoP for the initial training dataset, an RMS error of 0.03 was calculated for the training dataset.

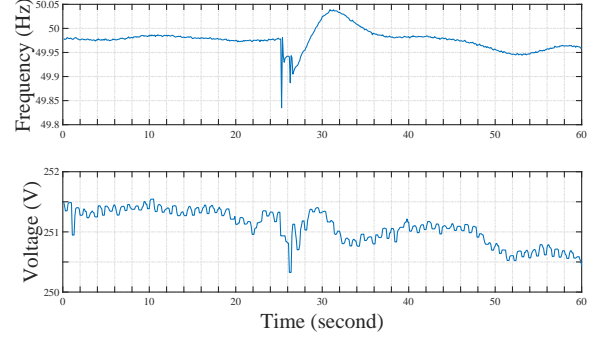
B. Power System Anomaly Detection Results

1) *Case 1 - Line Trip Event*: This case study presents a typical line trip event that was recorded on the PMU situated at Queen's, a plot of two of the chosen predictor variables (frequency and voltage) are shown in Fig 5 (a). It can be seen from this figure, that up until $t \approx 25$ seconds both system frequency and voltage are varying (as expected) around their nominal values, at $t \approx 25$ seconds the system frequency has a sharp decrease from 49.98 Hz to 49.84 Hz over 0.1 seconds. The frequency then returns to nominal again before experiencing an additional fall 49.89 Hz over 0.9 seconds before rising again and returning to nominal frequency around $t = 42$ seconds. Similar to frequency, voltage also experiences a sudden drop from 251.3 to 250.9 volts in 0.1 seconds, before further voltage drop to 250.3 volts is registered at the PMU, agreeing with the frequency plot previously analyzed.

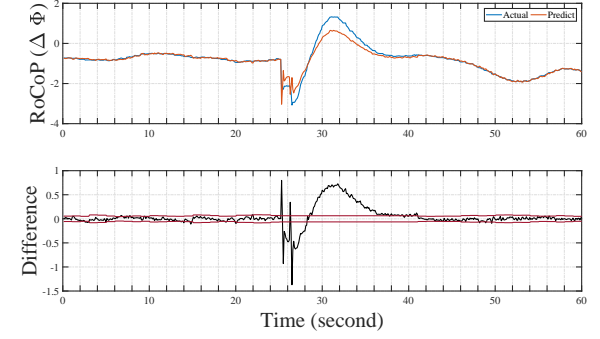
Monitoring results from the observed and predicted RoCoP are displayed in Fig 5 (b) which shows that pre-event there are small differences between observed and predicted RoCoP value, however these fall below the set threshold for anomaly detection in the PMU data. Agreeing with the variable plots in Fig 5, the method detects anomalous data when thresholds are violated at $t \approx 25$ seconds when the difference change from 0.1 to -0.9. Also from the monitoring results in Fig 5 (b) it can be observed that whilst the event is occurring, anomaly thresholds are held constant as to avoid contamination of the model by using event data in its construction. This allows a proper return to nominal value for the variables to be calculated at $t = 42$ seconds, again agreeing with the variable plot in 5. It should be noted that in the voltage variable plot the measurement from $t = 50 - 60$ seconds is much lower than some of the event voltage, however as can be seen in the result plot this does not cause the triggering of any anomaly.

2) *Case 2 - MITM Attack on Frequency Measurement*:

The modification of system frequency was implemented using the MAC-spoofing MITM technique described in Section III.



(a) Power System Variables



(b) Anomaly Detection monitoring

Fig. 5: Case 1: (a) frequency and voltage plots for line trip event, and (b) comparison of observed and predicted RoCoP values (top plot), and anomaly detection results showing the difference between observed and predicted values and thresholds for detection (bottom plot)

The attacker script is first set to the desired ramp and test period. The attacker PC is then plugged into the same physical network as the sending (or receiving) PMU - this would emulate a physical security breach in a substation or where the data is received. Once the script is executed it detects PMU data streams (using the common TCP port used for PMU data and the known footprint, or signature of C37.118.2 data), the system then initiates a MITM attack through spoofing operations, and places itself as a 'pathway' in the data stream. Once the path is established, the attacker can then initiate the manipulation of the data, taking the incoming data points, one at a time, and modifying them by the correct amount to replicate attack parameters. The modified packet is then sent to the intended recipient, where there are no indications that any data modification has taken place. This attack action can be applied to any of the system variables and over an indeterminate amount of time.

Displayed in Fig 6 (a) is a MITM attack that targeted live frequency data being reported by a PMU. In this case the frequency was increased from 49.5 Hz to 55 Hz over the period of 5 seconds. This type of frequency ramp could be interpreted as an islanding event (when a distributed generator continues to energize local loads after isolation from the main system), where generation is higher than load. If distributed generation assets were using this data for islanding protection, then they might disconnect or the system operator might intentionally isolate that part of the network.

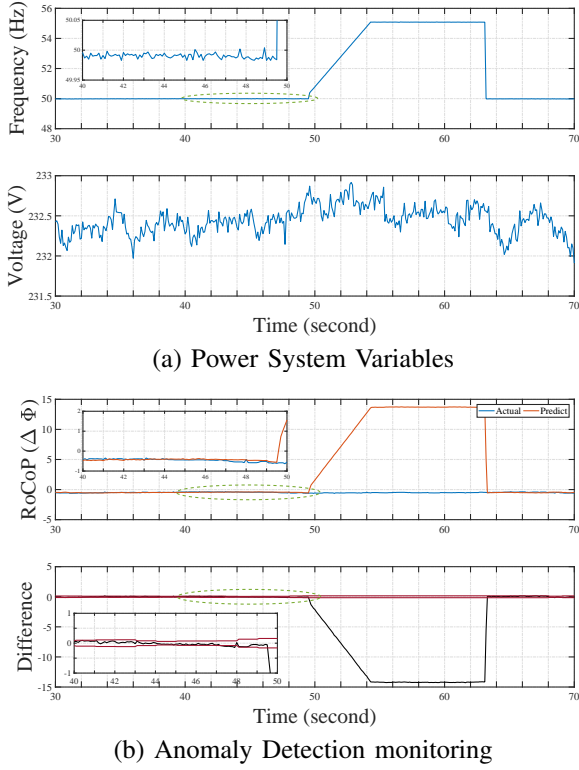


Fig. 6: Case 2: (a) frequency and voltage plots for frequency intrusion with insert of pre-anomaly system frequency also shown, and (b) anomaly detection results for frequency intrusion with inserts showing pre-anomaly monitoring

Anomaly detection monitoring results are presented in Fig 6 (b), the insert plots show results pre-anomaly and it can be seen that there is a small difference between observed and predicted RoCoP values (top plot), and as expected these fall between RMS error thresholds previously selected during the initialization phase of the methodology (bottom plot). It can be observed from the lower plot in Fig 6 (b) that at $t = 50$ seconds the difference between observed and predicted RoCoP values begin to exceed the set RMS error threshold and once a genuine anomaly has been detected in the PMU data, thresholds are held again. At $t = 63$ seconds the methodology detects the end of the anomaly, which corresponds with the system frequency plot in 6 (a). Once the anomalous data has been cleared from the system, anomaly thresholds begin to update again and difference between observed and predicted RoCoP falls between calculated thresholds.

IV. CONCLUSIONS AND FUTURE WORK

In this paper a methodology for the detection of local power system anomalies is proposed. The multiple linear regression method employs frequency, phase angle and voltage magnitude; collected from PMUs located across the Irish power system. Two case studies are presented as preliminary evaluations of the methodology, they demonstrate the potential for real-time applications and historical analysis. Setting a threshold of $3 \times$ the nominal RMS error was demonstrated as a reliable method of detecting anomalies in the PMU data.

The method proposed has several advantages including self calibration to local conditions, both locational and temporal, through the application of the RMS error. The quick ‘flagging mechanism’ has applications for real-time PMU data that could be applied at the central processing level or at remote locations. Flagging is used to highlight, anticipated and unanticipated, anomalies in data streams, before the relevant data is forwarded for more computationally intensive event categorization. The flagging mechanism could also be applied as a method of highlighting suspicious data. For example, if PMU data is being used as an input for a cyber-physical system, then control system may require two or more PMU streams to corroborate before actuating any physical system components.

This method is still under development and some of the conditions for anomaly detection and clearing are at present ad hoc; such as using a 0.5 second time delay before triggering a genuine anomaly. In practice, this may be slow and in some cases genuine anomalies may occur for less than this delay. While this delay is an arbitrary value, future investigation will determine a more optimal solution. At present this regression method has only been applied to RoCoP, the next step will be to apply it to other system variables. The method described can easily be applied to frequency, phase angle, voltage and their derivatives simultaneously; this may also give insight into the nature of the anomalies.

Future work will look at some enhancements to improve and further validate the methodology presented. Firstly, an investigation into optimal window size and frequency of regression model retraining will be carried out. Secondly, the simultaneous regression technique on multiple system variables will be carried out to identify if it can yield a more robust detection of anomalies, anomaly categorization and identify MITM attacks. Finally, the results from multiple locations will be combined to develop this from a local to a wide-area detection method.

ACKNOWLEDGMENT

The research is supported by a British Council Newton Institutional Links Programme grant with Helwan University, Egypt.

REFERENCES

- [1] “IEEE standard for synchrophasor data transfer for power systems,” Dec. 28 2011, IEEE C37.118.2-2011 (Revision of IEEE Std C37.118-2005) doi: 10.1109/IEEESTD.2011.6111222.
- [2] R. Khan, K. McLaughlin, D. Laverty, and S. Sezer, “Analysis of iec c37.118 and iec 61850-90-5 synchrophasor communication frameworks,” in *Power and Energy Society General Meeting (PESGM), 2016*. IEEE, 2016, pp. 1–5.
- [3] T. M. Mitchell *et al.*, “Machine learning. wcb,” 1997.
- [4] S. Chatterjee and A. S. Hadi, *Regression analysis by example*. John Wiley & Sons, 2015.
- [5] R. Akerkar, *Big data computing*. CRC Press, 2013.
- [6] D. Ruan, G. Chen, E. E. Kerre, and G. Wets, *Intelligent data mining: techniques and applications*. Springer Science & Business Media, 2005, vol. 5.
- [7] D. M. Laverty, R. J. Best, P. Brogan, I. Al Khatib, L. Vanfretti, and D. J. Morrow, “The OpenPMU platform for open-source phasor measurements,” *IEEE Trans. on Instrumentation and Measurement*, vol. 62, no. 4, pp. 701–709, 2013.

- [8] X. Zhao, D. M. Lavery, A. McKernan, D. J. Morrow, K. McLaughlin, and S. Sezer, "Gps-disciplined analog-to-digital converter for phasor measurement applications," *IEEE Transactions on Instrumentation and Measurement*, 2017.